



# EP3260: Machine Learning Over Networks

## Lecture 5: Centralized Nonconvex ML

Hossein S. Ghadikolaei

Division of Network and Systems Engineering  
School of Electrical Engineering and Computer Science  
KTH Royal Institute of Technology, Stockholm, Sweden

<https://sites.google.com/view/mlons2023/home>

February 2023

# Learning outcomes

- Recap of stochastic iterative algorithms for convex optimization
- Hardness of nonconvex problems
- Finding stationary points for (non-)smooth problems
- Escaping saddle points
- Finding an approximate local minima
- Convergence results for large-scale nonconvex optimization

# Outline

1. Nonconvex optimization
2. Finding stationary points for structured nonconvex problems
3. Finding stationary points for generic nonconvex problems
4. Finding a local minima
5. Supplements

# Recap of convex solvers

Our main optimization problem: minimize  $\frac{1}{N} \sum_{i \in [N]} f_i(\mathbf{w}) + r(\mathbf{w})$

Existence of global optimality and efficient solvers

## Deterministic solution algorithms

- gradient Oracle's load is **linear** with  $N$

- GD family for smooth problems

- subgradient and proximal methods for non-smooth (or composite) functions

## Stochastic solution algorithms

- gradient Oracle's load is **independent** of  $N$

- SGD family for smooth problems

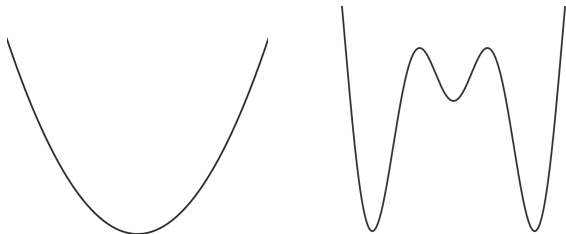
- noise reduction techniques and adaptive mini-batches, SVRG, and SAGA

- proximal methods for non-smooth (or composite) functions

# Outline

1. Nonconvex optimization
2. Finding stationary points for structured nonconvex problems
3. Finding stationary points for generic nonconvex problems
4. Finding a local minima
5. Supplements

# Nonconvexity



Nonconvex optimization can encode most problems

Local optimality **may not** necessarily imply global optimality

Proper initialization is very important in nonconvex optimization

First-order criteria ( $\|\nabla f(\mathbf{w})\|_2 \rightarrow 0$ )

necessary and sufficient conditions for convex

**only necessary condition** for nonconvex

# Nonconvex optimization

Subset-sum problem: find a subset of  $\{a_1, a_2, \dots, a_N\}$  that sums to  $b$ .  
Let's formulate it as a nonconvex optimization problem:

$$\underset{\mathbf{w} \in \{0,1\}^N}{\text{minimize}} \quad \left(\mathbf{a}^T \mathbf{w} - b\right)^2 + \mathbf{w}^T (\mathbf{1} - \mathbf{w})$$

Can we achieve the global minima (e.g., 0)?

It is NP-complete

OK! give up the global optimal point. Can we run GD to find a local minima?

Curse of dimensionality: finding a local optima becomes exponentially harder

⇒ Results in nonconvex optimization are not as strong as the convex case

# Nonconvex optimization

Subset-sum problem: find a subset of  $\{a_1, a_2, \dots, a_N\}$  that sums to  $b$ .  
Let's formulate it as a nonconvex optimization problem:

$$\underset{\mathbf{w} \in \{0,1\}^N}{\text{minimize}} \quad \left(\mathbf{a}^T \mathbf{w} - b\right)^2 + \mathbf{w}^T (\mathbf{1} - \mathbf{w})$$

Can we achieve the global minima (e.g., 0)?

It is NP-complete

OK! give up the global optimal point. Can we run GD to find a local minima?

Curse of dimensionality: finding a local optima becomes exponentially harder

⇒ Results in nonconvex optimization are not as strong as the convex case



# Nonconvex optimization

Subset-sum problem: find a subset of  $\{a_1, a_2, \dots, a_N\}$  that sums to  $b$ .  
Let's formulate it as a nonconvex optimization problem:

$$\underset{\mathbf{w} \in \{0,1\}^N}{\text{minimize}} \quad \left(\mathbf{a}^T \mathbf{w} - b\right)^2 + \mathbf{w}^T (\mathbf{1} - \mathbf{w})$$

Can we achieve the global minima (e.g., 0)?

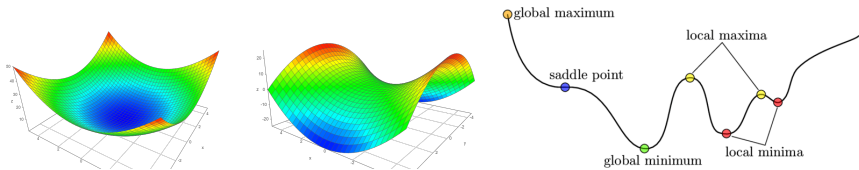
It is NP-complete

OK! give up the global optimal point. Can we run GD to find a local minima?

Curse of dimensionality: finding a local optima becomes exponentially harder

⇒ Results in nonconvex optimization are not as strong as the convex case

# Some definitions



**Stationary/critical points:**  $\{w \mid \nabla f(w) = 0\}$

**Local minima:** a critical point where  $\nabla^2 f(w) > 0$

**Local maxima:** a critical point where  $\nabla^2 f(w) < 0$

**Non-degenerate saddle points:** a critical point where  $\nabla^2 f(w)$  has strictly positive and negative eigenvalues

visualize gradient flow around a non-degenerate saddle point

Other types of saddle points (e.g., monkey saddle), usually harder to identify and treat!

# More definitions

Problem: minimize  $f(\mathbf{w}) = \frac{1}{N} \sum_{i \in [N]} f_i(\mathbf{w})$   
 $\mathbf{w} \in \mathcal{W}$

**Convex land:** **optimality gap**  $\mathbb{E} [\|\nabla f(\mathbf{w}_k)\|_2^2] \rightarrow 0$

expectation w.r.t. potential randomness of the algorithm

**Nonconvex land:** **stationarity gap**  $\mathbb{E} [\|\nabla f(\mathbf{w}_k)\|_2^2] \rightarrow 0$

Second-order necessary (2oN) point:  $\|\nabla f(\mathbf{w}_k)\|_2^2 \rightarrow 0$  &  $\nabla^2 f(\mathbf{w}_k) \geq 0$

Approximate 2oN point:  $\|\nabla f(\mathbf{w}_k)\|_2^2 \leq \epsilon_g$  ,  $\nabla^2 f(\mathbf{w}_k) \geq -\epsilon_H \mathbf{I}$  for small positive  $\epsilon_g, \epsilon_H$

**Complexity measure:**

gradient evaluations: # calls to incremental first-order oracle with input  $(\mathbf{w}, i)$  and output  $(f_i(\mathbf{w}), \nabla f_i(\mathbf{w}))$

#Hessian-vector products

#alternations among subproblems

## Basic assumptions

$f$  is bounded below:  $f(\mathbf{w}) \geq f_{\inf}$  for all  $\mathbf{w} \in \mathcal{W}$

Gradient and Hessian are Lipschitz continuous

$$\begin{aligned}\|\nabla f(\mathbf{w}_2) - \nabla f(\mathbf{w}_1)\|_2 &\leq L_g \|\mathbf{w}_2 - \mathbf{w}_1\|_2 \\ \|\nabla^2 f(\mathbf{w}_2) - \nabla^2 f(\mathbf{w}_1)\|_2 &\leq L_H \|\mathbf{w}_2 - \mathbf{w}_1\|_2\end{aligned}$$

Quadratic and cubic upper-bounds from Taylor's theorem ( $\forall \mathbf{v}, \mathbf{w} \in \mathcal{W}$ )

$$\begin{aligned}f(\mathbf{w} + \mathbf{v}) &\leq f(\mathbf{w}) + \nabla f(\mathbf{w})^T \mathbf{v} + \frac{L_g}{2} \|\mathbf{v}\|_2^2 \\ f(\mathbf{w} + \mathbf{v}) &\leq f(\mathbf{w}) + \nabla f(\mathbf{w})^T \mathbf{v} + \frac{1}{2} \mathbf{v}^T \nabla^2 f(\mathbf{x}) \mathbf{v} + \frac{L_H}{6} \|\mathbf{v}\|_2^3\end{aligned}$$

# Outline

1. Nonconvex optimization
2. Finding stationary points for structured nonconvex problems
3. Finding stationary points for generic nonconvex problems
4. Finding a local minima
5. Supplements

## Reformulation to something similar

- Minimum eigenvalue of a symmetric matrix  $A$ : 
$$\underset{\|w\|_2=1}{\text{minimize}} \quad w^T A w$$

Nonconvex on the Euclidean space, but easily solvable on the Riemannian manifolds [Smith, 1994]

Run your favorite convex solver on the reformulated convex problem

- Quasi-convex problems, invex problems, etc.

- Usually ad-hoc, no generic approach, problem specific

[Smith, 1994] S.T. Smith, "Optimization techniques on Riemannian manifolds," Fields Institute Communications, 1994.

# Can we use a simple GD?

Convexity implies  $f(\mathbf{w}_k) - f(\mathbf{v}) \leq \nabla f(\mathbf{w}_k)^T(\mathbf{w}_k - \mathbf{v})$  for all  $\mathbf{v}$

What if we have that inequality **only for one point**:  $\mathbf{v} = \mathbf{w}^*$ ?

This condition holds for most functions in the vicinity of any *local minima*  $\mathbf{w}^*$

What if  $\mathbf{w}^*$  is global minima?

**Run**  $\mathbf{w}_{k+1} = \mathbf{w}_k - \frac{\nabla f(\mathbf{w}_k)}{L_g}$  **and let**  $f(\mathbf{w}_k) - f(\mathbf{w}^*) \leq \nabla f(\mathbf{w}_k)^T(\mathbf{w}_k - \mathbf{w}^*)$

$\mathcal{O}(N)$  gradient oracle calls per iteration

$\mathcal{O}(L_g/\epsilon_g)$  iterations for  $L_g$ -smooth functions, so  $\mathcal{O}(NL_g/\epsilon_g)$  calls (**proof?**)

So why not keep using GD even for nonconvex?

poor scalability with  $N$

convergence to a stationary point, not necessarily a local minima

# Can we use a simple GD?

Convexity implies  $f(\mathbf{w}_k) - f(\mathbf{v}) \leq \nabla f(\mathbf{w}_k)^T (\mathbf{w}_k - \mathbf{v})$  for all  $\mathbf{v}$

What if we have that inequality **only for one point**:  $\mathbf{v} = \mathbf{w}^*$ ?

This condition holds for most functions in the vicinity of any *local minima*  $\mathbf{w}^*$

What if  $\mathbf{w}^*$  is global minima?

**Run**  $\mathbf{w}_{k+1} = \mathbf{w}_k - \frac{\nabla f(\mathbf{w}_k)}{L_g}$  **and let**  $f(\mathbf{w}_k) - f(\mathbf{w}^*) \leq \nabla f(\mathbf{w}_k)^T (\mathbf{w}_k - \mathbf{w}^*)$

$\mathcal{O}(N)$  gradient oracle calls per iteration

$\mathcal{O}(L_g/\epsilon_g)$  iterations for  $L_g$ -smooth functions, so  $\mathcal{O}(NL_g/\epsilon_g)$  calls (**proof?**)

So why not keep using GD even for nonconvex?

poor scalability with  $N$

convergence to a stationary point, not necessarily a local minima



# Can we use a simple GD?

Convexity implies  $f(\mathbf{w}_k) - f(\mathbf{v}) \leq \nabla f(\mathbf{w}_k)^T(\mathbf{w}_k - \mathbf{v})$  for all  $\mathbf{v}$

What if we have that inequality **only for one point**:  $\mathbf{v} = \mathbf{w}^*$ ?

This condition holds for most functions in the vicinity of any *local minima*  $\mathbf{w}^*$

**What if  $\mathbf{w}^*$  is global minima?**

**Run**  $\mathbf{w}_{k+1} = \mathbf{w}_k - \frac{\nabla f(\mathbf{w}_k)}{L_g}$  **and let**  $f(\mathbf{w}_k) - f(\mathbf{w}^*) \leq \nabla f(\mathbf{w}_k)^T(\mathbf{w}_k - \mathbf{w}^*)$

$\mathcal{O}(N)$  gradient oracle calls per iteration

$\mathcal{O}(L_g/\epsilon_g)$  iterations for  $L_g$ -smooth functions, so  $\mathcal{O}(NL_g/\epsilon_g)$  calls (**proof?**)

So why not keep using GD even for nonconvex?

poor scalability with  $N$

convergence to a stationary point, not necessarily a local minima

# One-point convexity [Allen-Zhu, ICML, 2017]

**Convex**  $f$

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{v} \rangle \geq f(\mathbf{w}) - f(\mathbf{v})$$

**$\mu$ -strongly convex**

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle \geq \mu \|\mathbf{w} - \mathbf{w}^*\|_2^2$$

$$\|\nabla f(\mathbf{w})\|_2^2 \geq 2\mu (f(\mathbf{w}) - f(\mathbf{w}^*))$$

**$\mu$ -strongly convex and  $L_g$ -smooth**

$$\begin{aligned} \langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle &\geq \frac{\mu}{2} \|\mathbf{w} - \mathbf{w}^*\|_2^2 \\ &\quad + \frac{1}{2L_g} \|\nabla f(\mathbf{w})\|_2^2 \end{aligned}$$

What if

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle \geq \beta (f(\mathbf{w}) - f(\mathbf{w}^*))$$

What if

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle \geq \beta \|\mathbf{w} - \mathbf{w}^*\|_2^2$$

two-layer NN [Li-Yuan, 2017]

$$\|\nabla f(\mathbf{w})\|_2^2 \geq \beta (f(\mathbf{w}) - f(\mathbf{w}^*))$$

(Polyak-Lojasiewicz condition)

finite sum minimization [Reddi-Sra-Poczos-Smola, 2016]

What if

$$\begin{aligned} \langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle &\geq \frac{\beta}{2} \|\mathbf{w} - \mathbf{w}^*\|_2^2 \\ &\quad + \gamma \|\nabla f(\mathbf{w})\|_2^2 \end{aligned}$$

dictionary learning [Arora-Ge-Ma-Moitra, 2015]

# One-point convexity [Allen-Zhu, ICML, 2017]

**Convex**  $f$

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{v} \rangle \geq f(\mathbf{w}) - f(\mathbf{v})$$

**$\mu$ -strongly convex**

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle \geq \mu \|\mathbf{w} - \mathbf{w}^*\|_2^2$$

$$\|\nabla f(\mathbf{w})\|_2^2 \geq 2\mu (f(\mathbf{w}) - f(\mathbf{w}^*))$$

**$\mu$ -strongly convex and  $L_g$ -smooth**

$$\begin{aligned} \langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle &\geq \frac{\mu}{2} \|\mathbf{w} - \mathbf{w}^*\|_2^2 \\ &\quad + \frac{1}{2L_g} \|\nabla f(\mathbf{w})\|_2^2 \end{aligned}$$

What if

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle \geq \beta (f(\mathbf{w}) - f(\mathbf{w}^*))$$

What if

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle \geq \beta \|\mathbf{w} - \mathbf{w}^*\|_2^2$$

two-layer NN [Li-Yuan, 2017]

$$\|\nabla f(\mathbf{w})\|_2^2 \geq \beta (f(\mathbf{w}) - f(\mathbf{w}^*))$$

(Polyak-Lojasiewicz condition)

finite sum minimization [Reddi-Sra-Poczos-Smola, 2016]

What if

$$\begin{aligned} \langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle &\geq \frac{\beta}{2} \|\mathbf{w} - \mathbf{w}^*\|_2^2 \\ &\quad + \gamma \|\nabla f(\mathbf{w})\|_2^2 \end{aligned}$$

dictionary learning [Arora-Ge-Ma-Moitra, 2015]

# One-point convexity, how to solve?

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{v} \rangle \geq f(\mathbf{w}) - f(\mathbf{v})$$

GD/SGD converges in  $\mathcal{O}\left(\frac{1}{\epsilon^2}\right)$

GD converges in  $\mathcal{O}\left(\frac{1}{\epsilon}\right)$  for smooth  $f$

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle \geq \mu \|\mathbf{w} - \mathbf{w}^*\|_2^2$$

GD/SGD converges in  $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$

$$\|\nabla f(\mathbf{w})\|_2^2 \geq 2\mu (f(\mathbf{w}) - f(\mathbf{w}^*))$$

GD converges in  $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$  for smooth  $f$

$$\begin{aligned} \langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle &\geq \frac{\mu}{2} \|\mathbf{w} - \mathbf{w}^*\|_2^2 \\ &\quad + \frac{1}{2L_g} \|\nabla f(\mathbf{w})\|_2^2 \end{aligned}$$

GD converges in  $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle \geq \beta (f(\mathbf{w}) - f(\mathbf{w}^*))$$

GD/SGD converges in  $\mathcal{O}\left(\frac{1}{\epsilon^2}\right)$

GD converges in  $\mathcal{O}\left(\frac{1}{\epsilon}\right)$  for smooth  $f$

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle \geq \beta \|\mathbf{w} - \mathbf{w}^*\|_2^2$$

GD/SGD converges in  $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$

$$\|\nabla f(\mathbf{w})\|_2^2 \geq \beta (f(\mathbf{w}) - f(\mathbf{w}^*))$$

GD converges in  $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$  for smooth  $f$

$$\begin{aligned} \langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle &\geq \frac{\beta}{2} \|\mathbf{w} - \mathbf{w}^*\|_2^2 \\ &\quad + \gamma \|\nabla f(\mathbf{w})\|_2^2 \end{aligned}$$

GD converges in  $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$

# One-point convexity, how to solve?

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{v} \rangle \geq f(\mathbf{w}) - f(\mathbf{v})$$

GD/SGD converges in  $\mathcal{O}\left(\frac{1}{\epsilon^2}\right)$

GD converges in  $\mathcal{O}\left(\frac{1}{\epsilon}\right)$  for smooth  $f$

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle \geq \mu \|\mathbf{w} - \mathbf{w}^*\|_2^2$$

GD/SGD converges in  $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$

$$\|\nabla f(\mathbf{w})\|_2^2 \geq 2\mu (f(\mathbf{w}) - f(\mathbf{w}^*))$$

GD converges in  $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$  for smooth  $f$

$$\begin{aligned} \langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle &\geq \frac{\mu}{2} \|\mathbf{w} - \mathbf{w}^*\|_2^2 \\ &\quad + \frac{1}{2L_g} \|\nabla f(\mathbf{w})\|_2^2 \end{aligned}$$

GD converges in  $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle \geq \beta (f(\mathbf{w}) - f(\mathbf{w}^*))$$

GD/SGD converges in  $\mathcal{O}\left(\frac{1}{\epsilon^2}\right)$

GD converges in  $\mathcal{O}\left(\frac{1}{\epsilon}\right)$  for smooth  $f$

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle \geq \beta \|\mathbf{w} - \mathbf{w}^*\|_2^2$$

GD/SGD converges in  $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$

$$\|\nabla f(\mathbf{w})\|_2^2 \geq \beta (f(\mathbf{w}) - f(\mathbf{w}^*))$$

GD converges in  $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$  for smooth  $f$

$$\begin{aligned} \langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle &\geq \frac{\beta}{2} \|\mathbf{w} - \mathbf{w}^*\|_2^2 \\ &\quad + \gamma \|\nabla f(\mathbf{w})\|_2^2 \end{aligned}$$

GD converges in  $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$

# One-point convexity, how to solve?

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{v} \rangle \geq f(\mathbf{w}) - f(\mathbf{v})$$

GD/SGD converges in  $\mathcal{O}\left(\frac{1}{\epsilon^2}\right)$

GD converges in  $\mathcal{O}\left(\frac{1}{\epsilon}\right)$  for smooth  $f$

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle \geq \mu \|\mathbf{w} - \mathbf{w}^*\|_2^2$$

GD/SGD converges in  $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$

$$\|\nabla f(\mathbf{w})\|_2^2 \geq 2\mu (f(\mathbf{w}) - f(\mathbf{w}^*))$$

GD converges in  $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$  for smooth  $f$

$$\begin{aligned} \langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle &\geq \frac{\mu}{2} \|\mathbf{w} - \mathbf{w}^*\|_2^2 \\ &\quad + \frac{1}{2L_g} \|\nabla f(\mathbf{w})\|_2^2 \end{aligned}$$

GD converges in  $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle \geq \beta (f(\mathbf{w}) - f(\mathbf{w}^*))$$

GD/SGD converges in  $\mathcal{O}\left(\frac{1}{\epsilon^2}\right)$

GD converges in  $\mathcal{O}\left(\frac{1}{\epsilon}\right)$  for smooth  $f$

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle \geq \beta \|\mathbf{w} - \mathbf{w}^*\|_2^2$$

GD/SGD converges in  $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$

$$\|\nabla f(\mathbf{w})\|_2^2 \geq \beta (f(\mathbf{w}) - f(\mathbf{w}^*))$$

GD converges in  $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$  for smooth  $f$

$$\begin{aligned} \langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle &\geq \frac{\beta}{2} \|\mathbf{w} - \mathbf{w}^*\|_2^2 \\ &\quad + \gamma \|\nabla f(\mathbf{w})\|_2^2 \end{aligned}$$

GD converges in  $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$

# One-point convexity, how to solve?

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{v} \rangle \geq f(\mathbf{w}) - f(\mathbf{v})$$

GD/SGD converges in  $\mathcal{O}\left(\frac{1}{\epsilon^2}\right)$

GD converges in  $\mathcal{O}\left(\frac{1}{\epsilon}\right)$  for smooth  $f$

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle \geq \mu \|\mathbf{w} - \mathbf{w}^*\|_2^2$$

GD/SGD converges in  $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$

$$\|\nabla f(\mathbf{w})\|_2^2 \geq 2\mu (f(\mathbf{w}) - f(\mathbf{w}^*))$$

GD converges in  $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$  for smooth  $f$

$$\begin{aligned} \langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle &\geq \frac{\mu}{2} \|\mathbf{w} - \mathbf{w}^*\|_2^2 \\ &\quad + \frac{1}{2L_g} \|\nabla f(\mathbf{w})\|_2^2 \end{aligned}$$

GD converges in  $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle \geq \beta (f(\mathbf{w}) - f(\mathbf{w}^*))$$

GD/SGD converges in  $\mathcal{O}\left(\frac{1}{\epsilon^2}\right)$

GD converges in  $\mathcal{O}\left(\frac{1}{\epsilon}\right)$  for smooth  $f$

$$\langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle \geq \beta \|\mathbf{w} - \mathbf{w}^*\|_2^2$$

GD/SGD converges in  $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$

$$\|\nabla f(\mathbf{w})\|_2^2 \geq \beta (f(\mathbf{w}) - f(\mathbf{w}^*))$$

GD converges in  $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$  for smooth  $f$

$$\begin{aligned} \langle \nabla f(\mathbf{w}), \mathbf{w} - \mathbf{w}^* \rangle &\geq \frac{\beta}{2} \|\mathbf{w} - \mathbf{w}^*\|_2^2 \\ &\quad + \gamma \|\nabla f(\mathbf{w})\|_2^2 \end{aligned}$$

GD converges in  $\mathcal{O}\left(\log\left(\frac{1}{\epsilon}\right)\right)$

# Outline

1. Nonconvex optimization
2. Finding stationary points for structured nonconvex problems
3. Finding stationary points for generic nonconvex problems
4. Finding a local minima
5. Supplements



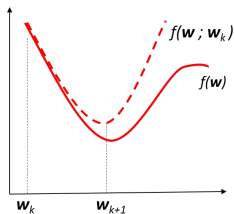
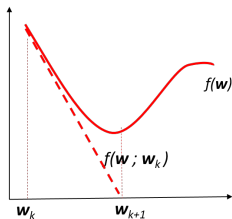
# 1. Successive approximation

Minimize an (approximate) surrogate fnc.  $\mathbf{w}_{k+1} = \arg \min_{\mathbf{w} \in \mathcal{W}} \tilde{f}(\mathbf{w}; \mathbf{w}_k)$ ,  
where  $\tilde{f}(\mathbf{w}; \mathbf{w}_k)$  is the approximation of  $f$  at  $\mathbf{w}_k$

- successive linear approximation:  $\tilde{f} = f(\mathbf{w}_k) + \nabla f(\mathbf{w}_k)^T (\mathbf{w} - \mathbf{w}_k)$
- successive quadratic approximation:

$$\tilde{f} = f(\mathbf{w}_k) + \nabla f(\mathbf{w}_k)^T (\mathbf{w} - \mathbf{w}_k) + \frac{1}{2}(\mathbf{w} - \mathbf{w}_k)^T \nabla^2 f(\mathbf{w}_k) (\mathbf{w} - \mathbf{w}_k)$$

-GD acts as successive quadratic upper-bound minimization when we replace  $\nabla^2 f(\mathbf{w}_k) \leq L\mathbf{I}$  (see Slide 2-24)

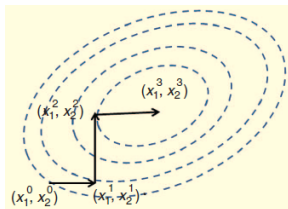


## 2. Coordinate descent (alternating methods)

Minimize only a block of coordinates while fixing others

$$\begin{aligned}\mathbf{w}_{(i)}^{k+1} &:= \arg \min_{\mathbf{w}_{(i)}} f(\mathbf{w}_{(1)}^k, \dots, \mathbf{w}_{(i-1)}^k, \mathbf{w}_{(i)}, \mathbf{w}_{(i+1)}^k, \dots, \mathbf{w}_{(d)}^k) \\ &= \arg \min_{\mathbf{w}_{(i)}} f(\mathbf{w}_{(i)}; \mathbf{w}^k)\end{aligned}$$

where  $\mathbf{w}_{(i)}$  is  $i$ -th block of coordinates



[Hong-Razaviyayn-Luo-Pang, 2016]

Cyclic update rules, random coordinate selection, etc.

Convergence rate of  $\mathcal{O}(dL_g/\epsilon_g)$  when each coordinate is  $L_g$ -smooth, namely  $[\nabla^2 f(\mathbf{w})]_{ii} \in [0, L_g]$ , convergence rate of  $\mathcal{O}(dL_g/\mu \log \epsilon_g^{-1})$  for  $\mu$ -strongly convex, [ShalevShwartz-Zhang, 2012]

What if the function is nonconvex/non-smooth in some coordinates?

### 3. Block successive upper-bound minimization

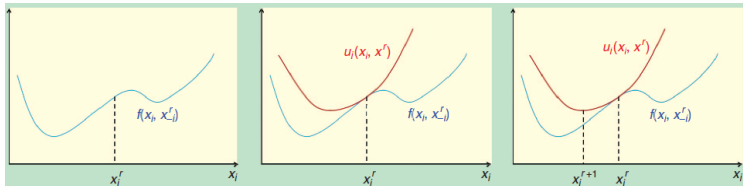
Combining block coordinate descent and successive convex approximation to handle non-convexity in  $f(\mathbf{w}_{(i)}; \mathbf{w}^k)$

$$\mathbf{w}_{(i)}^{k+1} = \arg \min_{\mathbf{w}_{(i)}} \tilde{f}(\mathbf{w}_{(i)}; \mathbf{w}^k)$$

where  $\tilde{f}$  is a **strongly convex smooth upperbound** of  $f$  at  $\mathbf{w}^k$ .

Example: proximal upper bound  $\tilde{f}(\mathbf{w}_{(i)}; \mathbf{w}^k) := f(\mathbf{w}_{(i)}; \mathbf{w}^k) + \frac{\gamma}{2} \|\mathbf{w}_{(i)} - \mathbf{w}_{(i)}^k\|_2^2$

Straightforward extension to composite functions ( $f = g + h$  for smooth  $g$  and non-smooth  $h$  through proximal mapping)



[Hong-Razaviyayn-Luo-Pang, 2016]

# Coordinate descent vs SGD

## Coordinate Descent

Pick coordinates  $\zeta_k$  and update  $\mathbf{w}_{k+1} = \mathbf{w}_k$  for coordinates  $i \neq \zeta_k$  and

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla_i f(\mathbf{w}_k)$$

for coordinates  $i = \zeta_k$

Compute all gradients for coordinates  $\zeta_k$

Guaranteed improvement in every iteration, consequently (usually) easier design and analysis

## Stochastic Gradient Descent

Pick coordinate(s)  $\zeta_k$  and update

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla f_{\zeta_k}(\mathbf{w}_k)$$

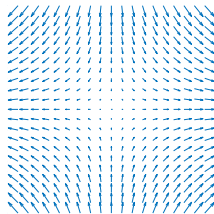
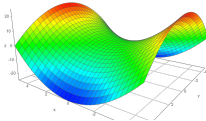
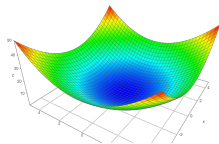
Compute one gradient  $\nabla f_{\zeta_k}$  for all coordinates

Not robust in general, need variance reduction techniques (like SVRG) to stabilize SGD

# Outline

1. Nonconvex optimization
2. Finding stationary points for structured nonconvex problems
3. Finding stationary points for generic nonconvex problems
4. Finding a local minima
5. Supplements

# How to escape non-degenerate saddle points



How to find an approximate 2oN point ( $\|f(\mathbf{w}_k)\|_2^2 \leq \epsilon_g, \nabla^2 f(\mathbf{w}_k) \geq -\epsilon_H \mathbf{I}$ )?

## Hessian-free approaches

1) Perturbed GD:  $\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k) + \text{noise}$

$\log(d)$  dependency on the parameter size [Jin-Ge-Netrapalli-Kakade-Jordan, 2017]

2) SGD

## Hessian-based approaches

Faster reaction to saddle points using curvature information, expensive iterations could be very efficient [Allen-Zhu, 2018]

# A generic Hessian-based algorithm

1. If  $\|\nabla f(\mathbf{w}_k)\|_2 > \epsilon_g$ , run your favorite algorithm (say GD with step-size  $1/L_g$  to get closer to a stationary point
2. Otherwise, if  $\nabla^2 f(\mathbf{w}_k) \not\preceq \epsilon_h \mathbf{I}$ , find the eigenvector of the most negative eigenvalue ( $\lambda_{\min}^k$ ) of  $\nabla^2 f(\mathbf{w}_k)$ , namely

$$\|\mathbf{v}_k\|_2^2 = 1, \quad \mathbf{v}^T \nabla^2 f(\mathbf{w}_k) \mathbf{v} = \lambda_{\min}^k, \quad \nabla f(\mathbf{w}_k)^T \mathbf{v} \leq 0$$

3. Move toward that direction (why?), e.g., by

$$\mathbf{w}_{k+1} = \mathbf{w}_k + \frac{2|\lambda_{\min}^k|}{L_H} \mathbf{v}$$

4. Otherwise, terminate.

⇒ The number of iterations is at most  $\max\left(\frac{2L_g}{\epsilon_g^2}, \frac{3L_H^2}{2\epsilon_h^3}\right) (f(\mathbf{w}_0) - f_{\inf})$

Proof: see the board

<sup>1</sup>Notice from Taylor expansion that  $f(\mathbf{w}_k + \mathbf{v}) \approx f(\mathbf{w}_k) + \nabla f(\mathbf{w}_k)^T \mathbf{v} + \frac{1}{2} \mathbf{v}^T \nabla^2 f(\mathbf{w}_k) \mathbf{v}$ . Step 2 finds a direction ( $\mathbf{v}$ ) that gives the highest reduction to  $f$ . Do we need to really find the minimum eigenvalue? or a strong negative one is enough?

# A generic Hessian-based algorithm

1. If  $\|\nabla f(\mathbf{w}_k)\|_2 > \epsilon_g$ , run your favorite algorithm (say GD with step-size  $1/L_g$  to get closer to a stationary point
2. Otherwise, if  $\nabla^2 f(\mathbf{w}_k) \not\preceq \epsilon_h \mathbf{I}$ , find the eigenvector of the most negative eigenvalue ( $\lambda_{\min}^k$ ) of  $\nabla^2 f(\mathbf{w}_k)$ , namely

$$\|\mathbf{v}_k\|_2^2 = 1, \quad \mathbf{v}^T \nabla^2 f(\mathbf{w}_k) \mathbf{v} = \lambda_{\min}^k, \quad \nabla f(\mathbf{w}_k)^T \mathbf{v} \leq 0$$

3. Move toward that direction (**why?**), e.g., by

$$\mathbf{w}_{k+1} = \mathbf{w}_k + \frac{2|\lambda_{\min}^k|}{L_H} \mathbf{v}$$

4. Otherwise, terminate.

$\Rightarrow$  The number of iterations is at most  $\max\left(\frac{2L_g}{\epsilon_g^2}, \frac{3L_H^2}{2\epsilon_h^3}\right) (f(\mathbf{w}_0) - f_{\inf})$

Proof: see the board

<sup>1</sup>Notice from Taylor expansion that  $f(\mathbf{w}_k + \mathbf{v}) \approx f(\mathbf{w}_k) + \nabla f(\mathbf{w}_k)^T \mathbf{v} + \frac{1}{2} \mathbf{v}^T \nabla^2 f(\mathbf{w}_k) \mathbf{v}$ . Step 2 finds a direction ( $\mathbf{v}$ ) that gives the highest reduction to  $f$ . Do we need to really find the minimum eigenvalue? or a strong negative one is enough?



# Foods for thought

1. Consider iterations of the acceleration methods developed for an  $L_g$ -smooth  $\mu$ -strongly convex problem. How to use them in a smooth nonconvex setting?

Estimate the minimum eigenvalue of  $\nabla^2 f(w_k)$ , add a proper convex quadratic term to the objective.

Find  $L_g$  and  $\mu$  for this new objective function and run Nesterov's iterations. Do we converge? at which rate? What about a non-smooth nonconvex objective?

2. Lanczos algorithm efficiently finds the smallest eigenvalue of a symmetric matrix  $A$ . Check the definition of Krylov subspace, generated by  $A$ . Find a vector in the Krylov subspace that minimizes  $z^T A z$ . Show that you can control the accuracy of the approximated eigenvalue by changing the degree of the Krylov subspace.

3. In our examples of successive convex approximation, we have used a *low-order* Taylor approximation around  $w_k$ . Can we trust this approximation anyway?

What about finding minimizer over a region to which we can trust? Define a trust region as a ball centered at  $w_k$  with radius  $\Delta_k$ .

Now, modify the successive quadratic approximation of slide 4-16 to find  $v$  in our trusted region. Adjust  $\Delta_k$  to ensure a sufficient decrease in  $f$  at each iteration

Congratulation! you have discovered the famous Trust Region method!

Extend to cubic regularization by replacing the 3rd term of the Taylor expansion by  $M_k \|v\|^3$  for some positive  $M_k$  (Nesterov and Polyak, 2006)? Observe that  $M_k \geq L_h/6$  leads to a successive upper-bound minimization.

4. Check <http://www.offconvex.org> and <https://www.facebook.com/nonconvex>

# Which algorithm to choose?

## CA3: Training a deep neural network

Consider optimization problem

$$\underset{\mathbf{W}_1, \mathbf{W}_2, \mathbf{w}_3}{\text{minimize}} \quad \frac{1}{N} \sum_{i \in [N]} \|\mathbf{w}_3 s(\mathbf{W}_2 s(\mathbf{W}_1 \mathbf{x}_i) - \mathbf{y}_i)\|_2^2,$$

where  $s(x) = 1/(1 + \exp(-x))$ . You may add your choice of regularizer.

Consider both “Individual household electric power consumption” and “Greenhouse gas observing network” datasets.

- 1) Try to solve this optimization task with proper choices of size of decision variables (matrix  $\mathbf{W}_1$ , matrix  $\mathbf{W}_2$ , and vector  $\mathbf{w}_3$ ) using GD, perturbed GD, SGD, SVRG, and block coordinate descent. For the SGD method, you may use the mini-batch version.
- 2) Compare these solvers in terms complexity of hyper-parameter tuning, convergence time, convergence rate (in terms of # outer-loop iterations), and memory requirement

## Some references

- Stephen J. Wright, "Coordinate descent algorithm" , Math. Program., 2015.
- S. Reddi, S. Sra, B. Póczos, and A. Smola, "Proximal stochastic methods for nonsmooth nonconvex finite-sum optimization," , NIPS, 2016.
- S. Shalev-Shwartz and T. Zhang, "Stochastic dual coordinate ascent methods for regularized loss minimization," JMLR, 2012.
- Y. Li and Y. Yuan, "Convergence analysis of two-layer neural networks with ReLU activation," , NIPS, 2017.
- M. Razaviyayn, M. Hong, and Z.-Q. Luo, "A unified convergence analysis of block successive minimization methods for nonsmooth optimization" , SIAM J. Optim., 2013.
- M. Hong, M. Razaviyayn, Z. Q. Luo and J. S. Pang, "A unified algorithmic framework for block-structured optimization involving big data: With applications in machine learning and signal processing," IEEE SPM, 2016.
- Z. Allen-Zhu, "Recent advances in stochastic convex and non-convex optimization," ICML Tutorial, 2017.
- S. Arora, R. Ge, T. Ma, and A. Moitra, "Simple, efficient, and neural algorithms for sparse coding," JMLR, 2015.
- C. Jin, R. Ge, P. Netrapalli, S.M. Kakade, and M.I. Jordan, "How to escape saddle points efficiently," ICML, 2017.
- Z. Allen-Zhu, "Natasha 2: Faster non-convex optimization than SGD," NIPS, 2018.
- Y. Zhang, P. Liang, and M. Charikar, "A hitting time analysis of stochastic gradient Langevin dynamics," COLT, 217.

# Outline

1. Nonconvex optimization
2. Finding stationary points for structured nonconvex problems
3. Finding stationary points for generic nonconvex problems
4. Finding a local minima
5. Supplements

## Application: Matrix factorization

*Matrix Factorization:* Given a data matrix  $\mathbf{D} \in \mathbb{R}^{N \times M}$ , find low-rank matrices,  $\mathbf{X} \in \mathbb{R}^{N \times k}$  and  $\mathbf{Y} \in \mathbb{R}^{M \times k}$ , such that  $\mathbf{D} \approx \mathbf{X}\mathbf{Y}^T$

$$\min_{\mathbf{X}, \mathbf{Y}} \|\mathbf{D} - \mathbf{X}\mathbf{Y}^T\|_F^2$$

Solved using BCD.

- Given  $\mathbf{Y}_k$  update  $\mathbf{X}_{k+1} := \min_{\mathbf{X}} \|\mathbf{D} - \mathbf{X}\mathbf{Y}_k^T\|_F^2$
- Given  $\mathbf{X}_{k+1}$  update  $\mathbf{Y}_{k+1} := \min_{\mathbf{Y}} \|\mathbf{D} - \mathbf{X}_{k+1}\mathbf{Y}^T\|_F^2$
- subproblems are strongly convex: convergence follows from BCD
- a.k.a. Alternating Least Squares
- BCD is basis of many algorithm in recommendation system: non-negative MF, sparse MF, etc.

## Application: Sparse dictionary learning

Given a data matrix  $\mathbf{D} \in \mathbb{R}^{N \times M}$ , find a dictionary  $\mathbf{X}\mathbf{Y}^T$ , that sparsely represents the data matrix,

$$\min_{\mathbf{X}, \mathbf{Y}} \|\mathbf{D} - \mathbf{X}\mathbf{Y}^T\|_F^2 + \lambda \|\mathbf{X}\|_1 \quad \text{s. t.} \quad \|\mathbf{Y}\|_F \leq \beta$$

## Special cases

Many known algorithms are special cases of block successive upperbound minimization (BSUM)

*Difference of convex (DC) programming:*

$\arg \min_{\mathbf{w}} f(\mathbf{w}) = g_1(\mathbf{w}) - g_2(\mathbf{w})$ , where  $g_1$  and  $g_2$  convex

$\mathbf{w}^{k+1} = \arg \min_{\mathbf{w}} g_1(\mathbf{w}) - (\nabla g_2(\mathbf{w}^k)^T (\mathbf{w} - \mathbf{w}^k)) - g_2(\mathbf{w}^k)$

- *Convex Concave Procedure*

*Block coordinate descent (BCD):*

Select the upperbound in BSUM as the function:

$u_i(\mathbf{w}_i, \mathbf{w}_{-i}^k) = f(\mathbf{w}_i, \mathbf{w}_{-i}^k)$

- we recover BCD

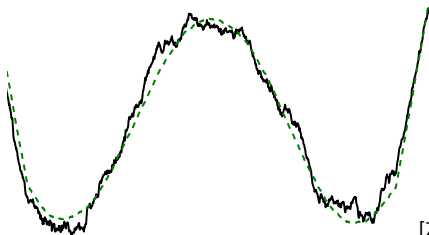
# Smooth and optimize

Consider a complicated non-smooth function  $f(\mathbf{w})$

One may use a smooth approximation  $g(\mathbf{w}) = \mathbb{E}_{\theta \sim \mathcal{B}(0, R)}[f(\mathbf{w} + \theta)]$  for some distribution of  $\theta$  defined on Borel set  $\mathcal{B}$

Random initialization on the obtained smooth function followed by SGD

This approach escapes suboptimal local minima that only exist in the empirical risk (black curve) not the population risk (green curve)



[Zhang-Liang-Charikar, 2017]